

# Improving ML training sets with Active Learning in



ML4Astro Catania 2024

Anais Möller  
ARC DECRA Fellow



**Australian Government**  
**Australian Research Council**

*Acknowledging the traditional owners of Swinburne's land,  
Wurundjeri People of the Kulin Nation*



# Rubin LSST

Now

## Raw Data

Sequential 30s image, 20TB/night

60s

## Prompt Data Product

Difference Image Analysis  
Alerts: up to 10 million per night

Public data!

24h

## Prompt Products DataBase

Images, Object and Source catalogs from DIA  
Orbit catalog for ~6 million Solar System bodies

Year

## Annual Data Release

Accessible via the LSST Science Platform &  
LSST Data Access Centers.

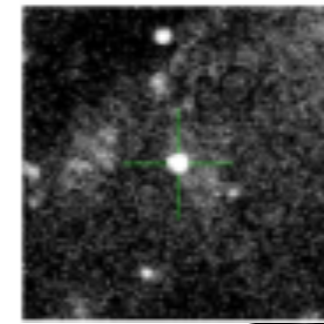
End

## Final 10yr Data Release

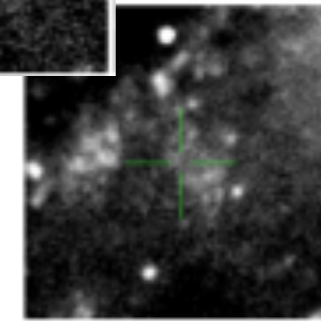
Images: 5.5 million x 3.2 Gpx  
Catalog: 15PB, 37 billion objects

LSST PIs and teams

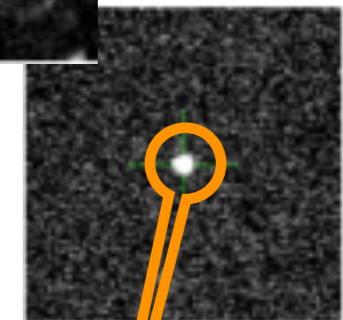
Credits: E. Bellm



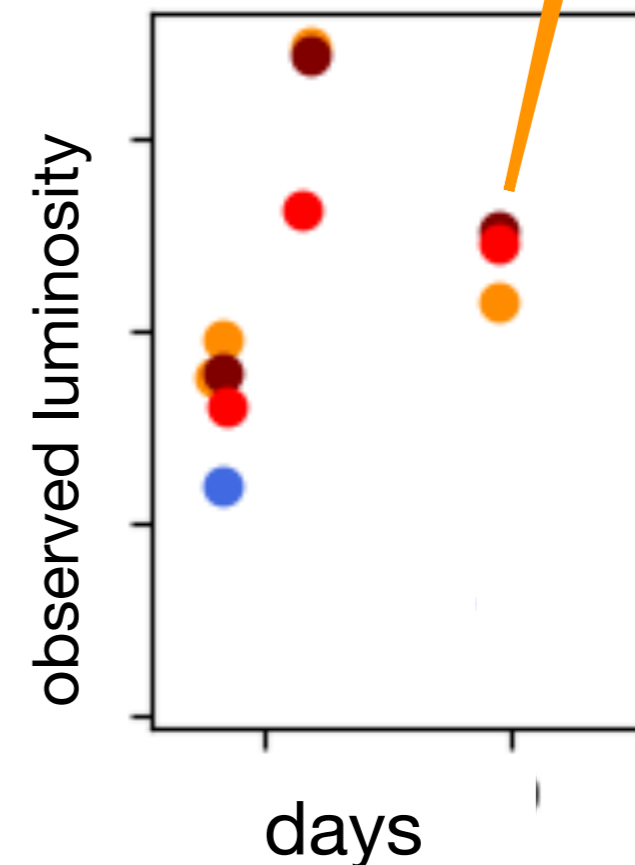
Observation



Template



Difference



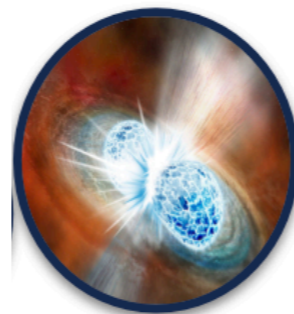
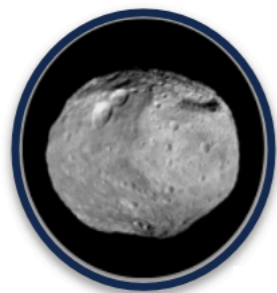
# Rubin LSST

10 million t-domain candidates/night

*Brokers such as*



*Hosted at CC-IN2P3 for Rubin*





LAPP, Ancey - 19 and 20 May 2022



2nd Fink Collaboration Meeting  
8 - 10 January 2024, IJCLab - France

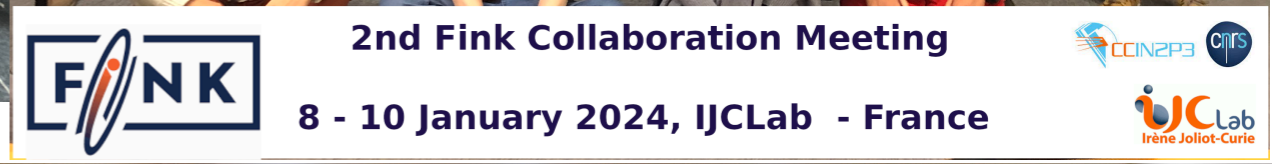


Photo: Carl Knox - OzGrav, Centre of Excellence for Gravitational Wave Discovery

3-5 May 2023, Melbourne - Australia



6-10 May 2024 - CBPF, Rio de Janeiro - Brazil

Fink-Brazil Workshop



*broad-science  
collaboration*

*community driven,  
open to everyone*

*PIs & management team  
Peloton, Ishida, AM*



Community driven, join us!



## Supernovae

SN Ia, SN, PISN

**Allam+2023, Leoni+2022, Möller+2022**

*CNRS MITI grant*

## Kilonovae

Orphan, GW+KN, Fast transients



**Grandma+Fink 2022,2023**

**Biswas+2022,2023**

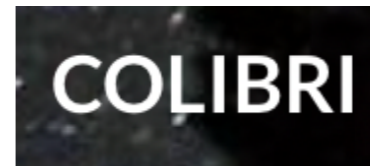
## Nuclear transients

AGN, TDE

**VAST** *Russeil+2022*

*Gondhalekar+ in prep.*

## Gamma Ray Bursts



## Solar System Objects

Discovery, tracking



*Le Montagner+ 2023*

*Carry+2024*

*CNRS MITI grant*

## Space Awareness

Satellite glints

*Karpov+2023,2022*

## Anomalies

**Pruzhinskaya+ in prep.**

## LSST simulations: ELAsTiCC

**Fraga+ 2024**

# Type Ia supernovae

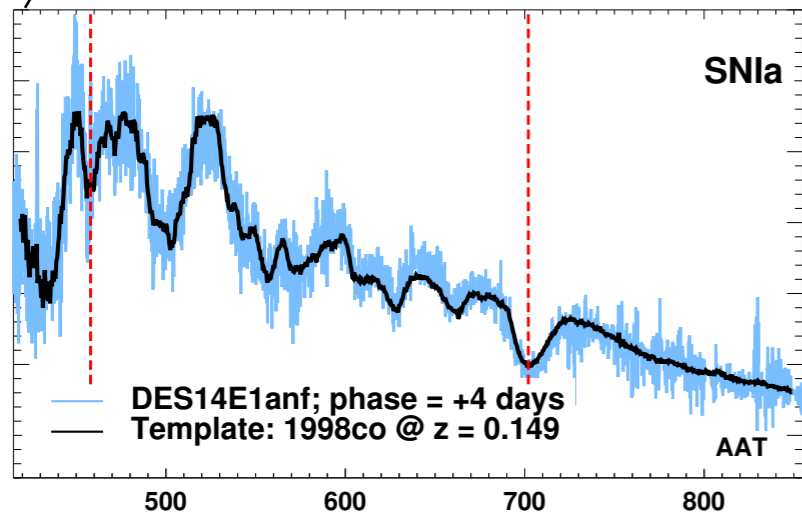
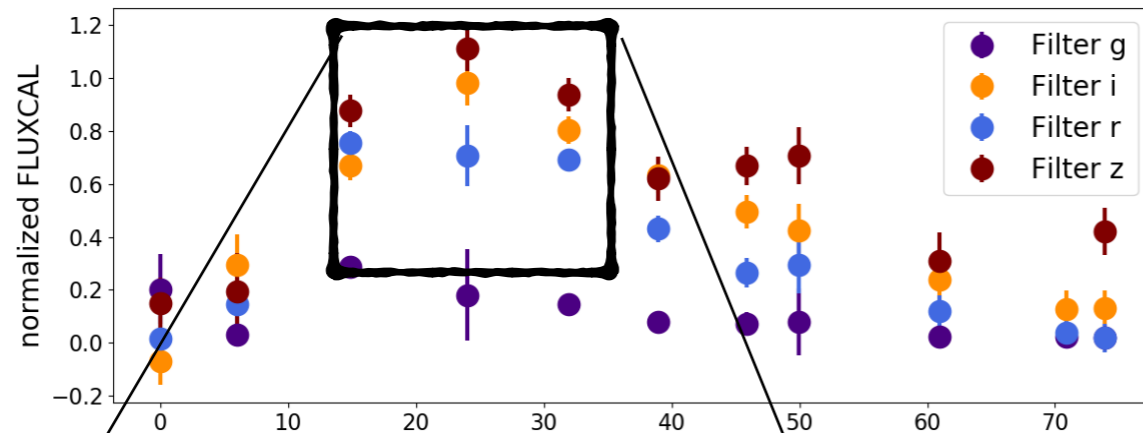


Homogeneous spectral and photometric properties

Direct measurement cosmic expansion

Rubin > 1 million

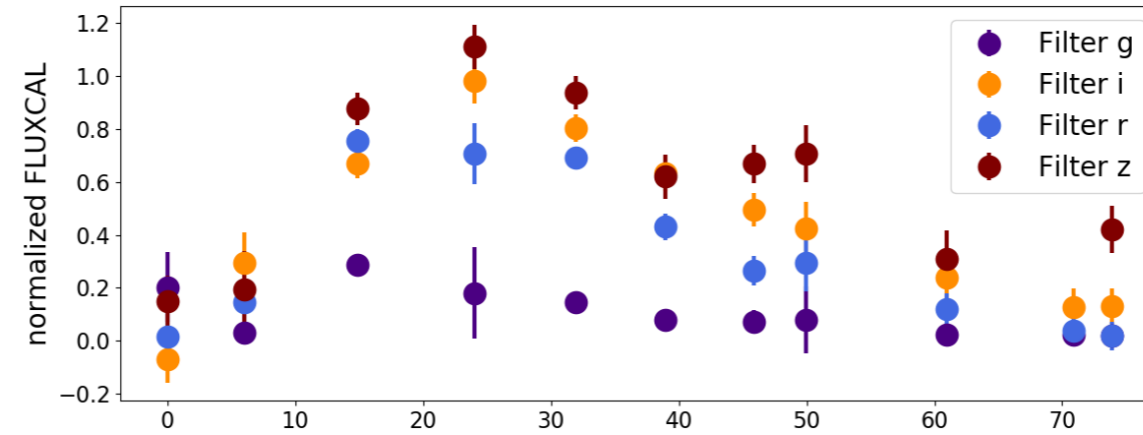
DES ~ 2k



**415 SNe Ia**

(expected 2,000 detected)

## Spectroscopic classification



ML algorithm



Photometric classification



1. Supervised + simulations
2. Unsupervised
3. AL for improving training sets

**1. Supervised + simulations**

2. Unsupervised

3. AL for improving training sets

# SuperNNova

Recurrent Neural Networks

Bayesian Neural Networks

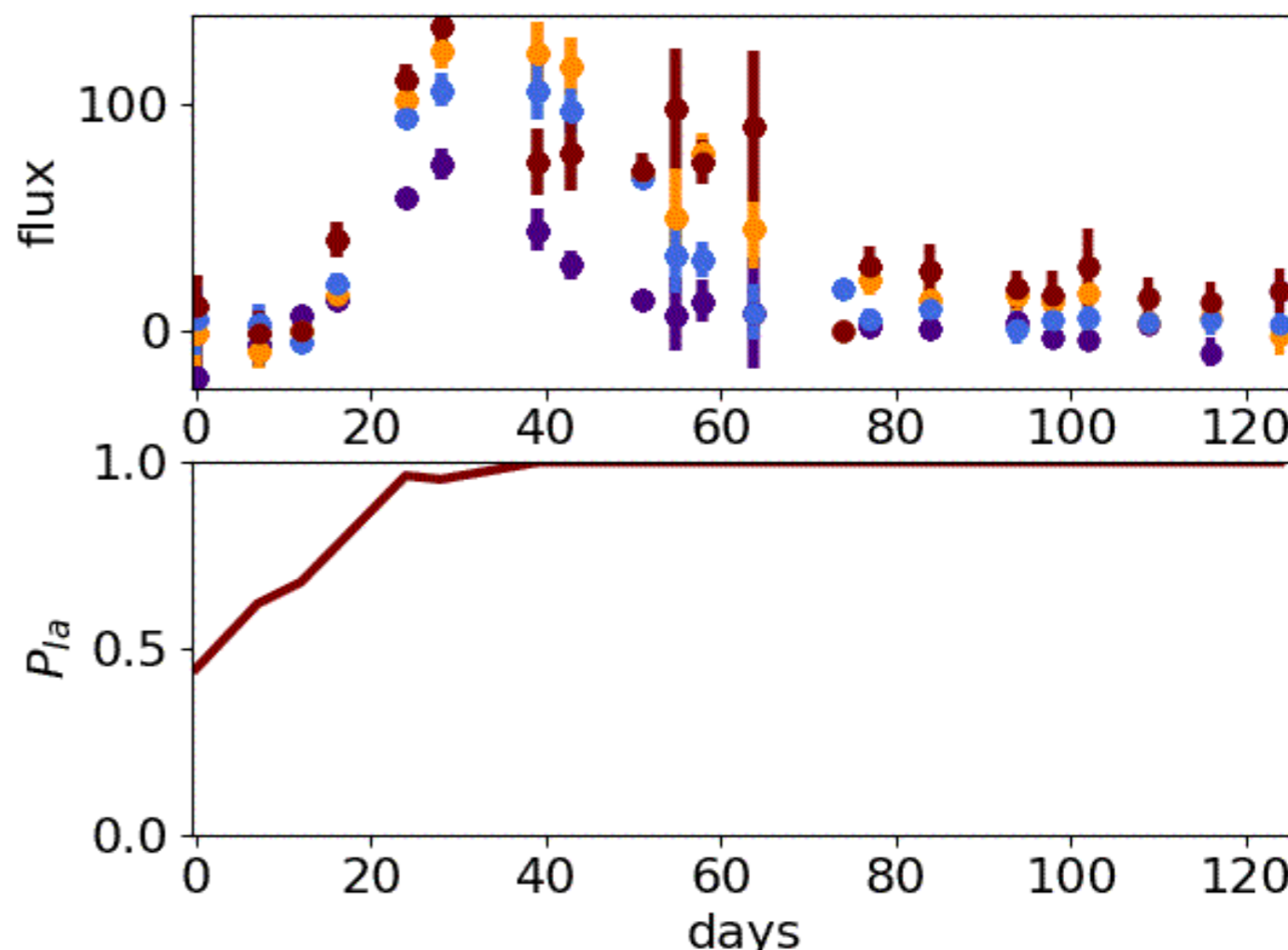
(MC dropout & Bayes by Backdrop)

*Möller+ 2019, 2022b*

Handles irregular time series

Used with DES, PS, ZTF data

Soon Rubin!

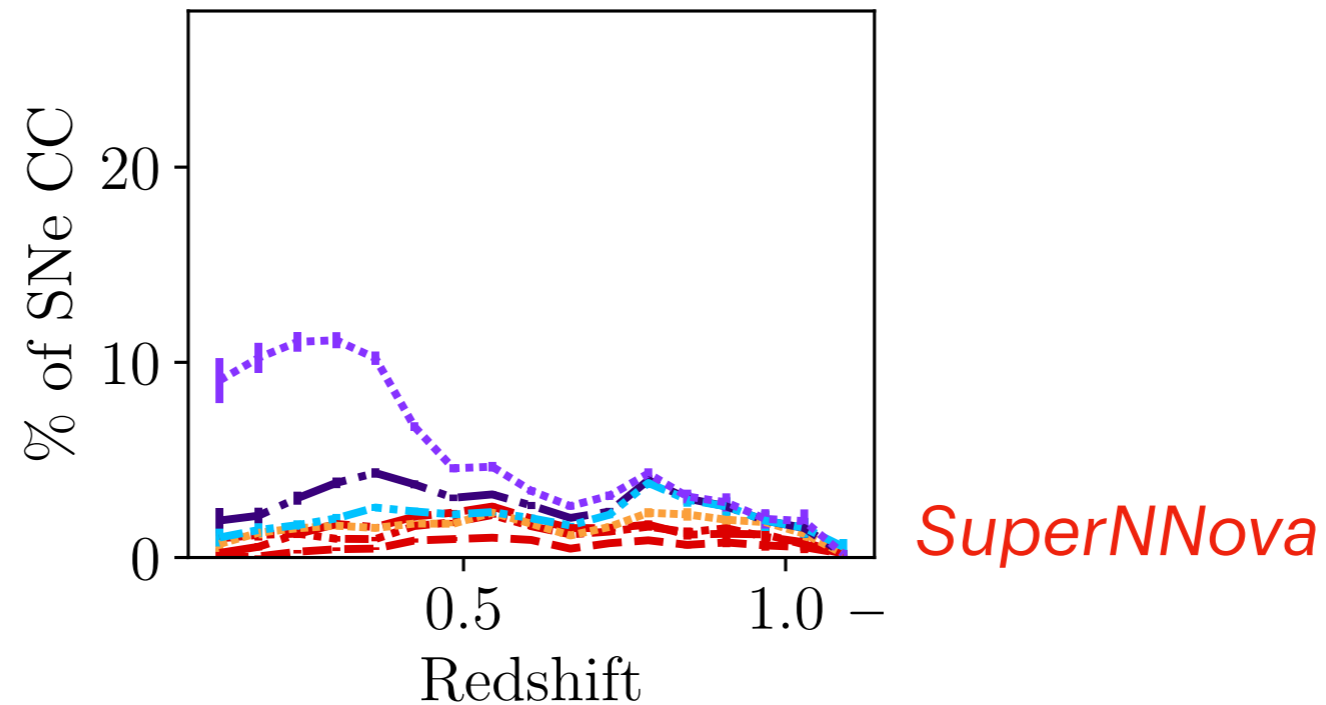


 *SuperNNova*

*Möller+2019*

*Pytorch, many configurations possible*

# DES: Light-curves + host-galaxy redshifts



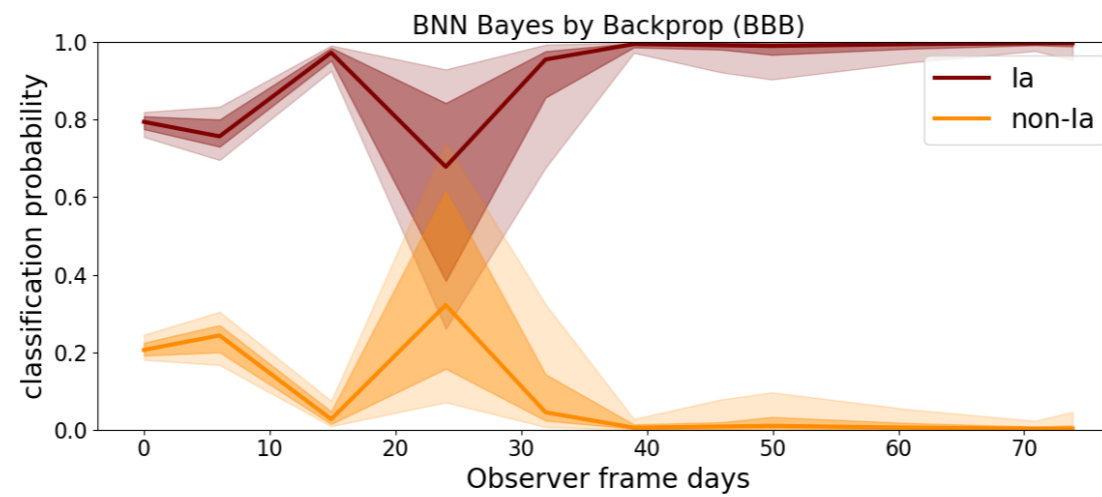
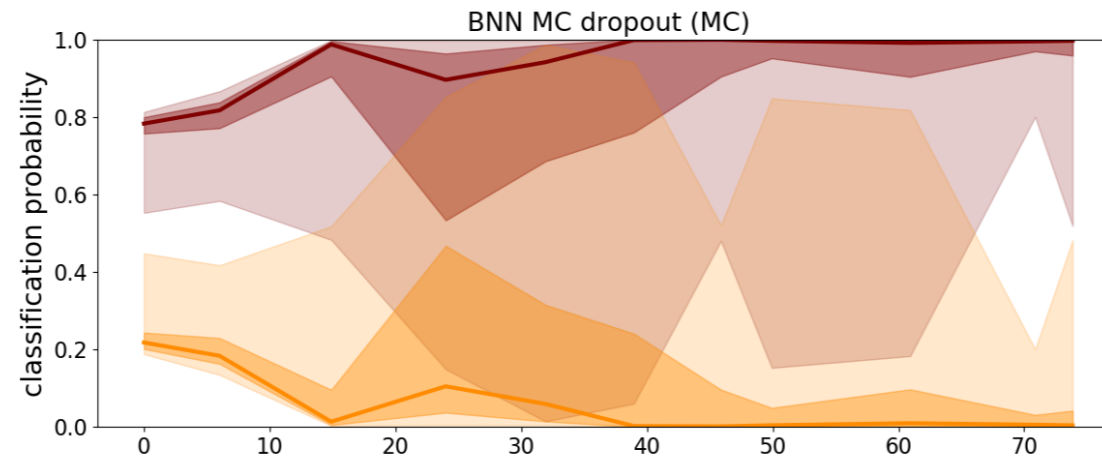
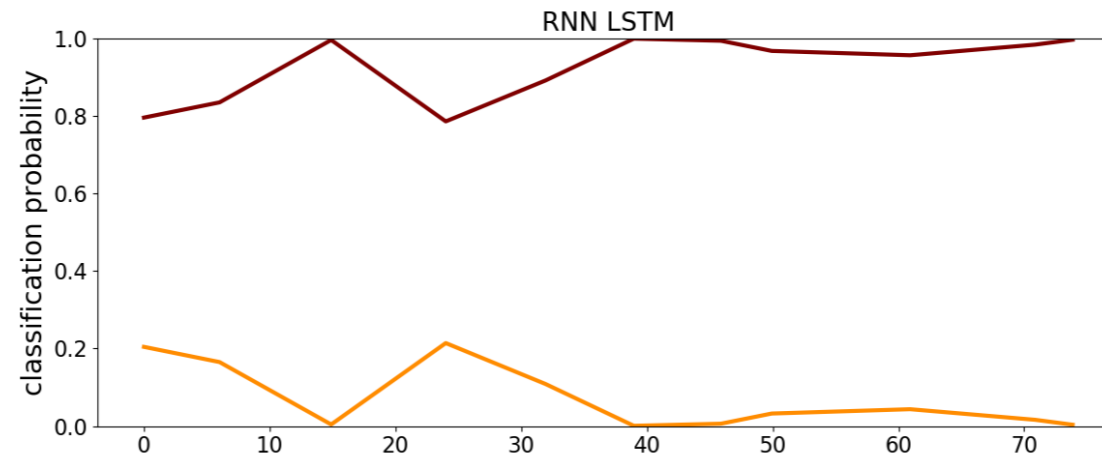
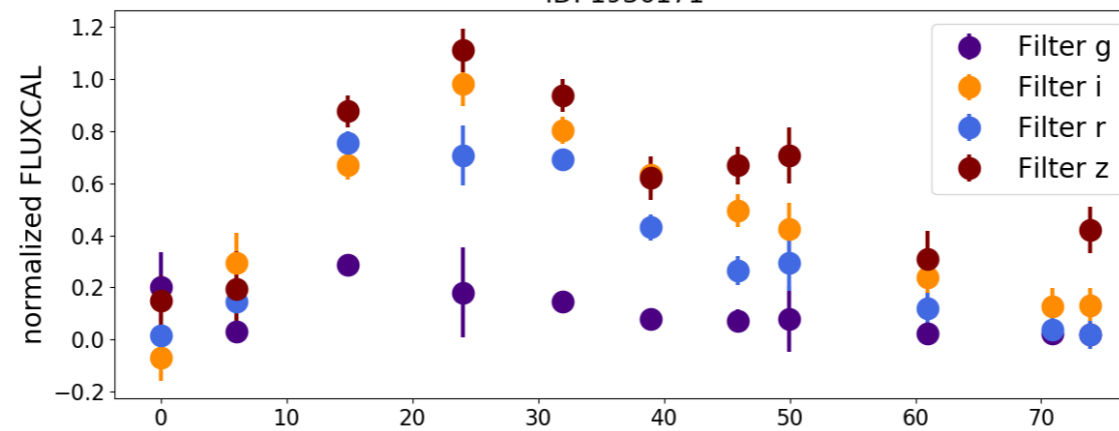
*Vincenzi, Sullivan, Möller et al. 2021*

State-of-the-art simulations

**>98% accuracy**

*Vincenzi+2020*

ID: 1936171



Möller+2022

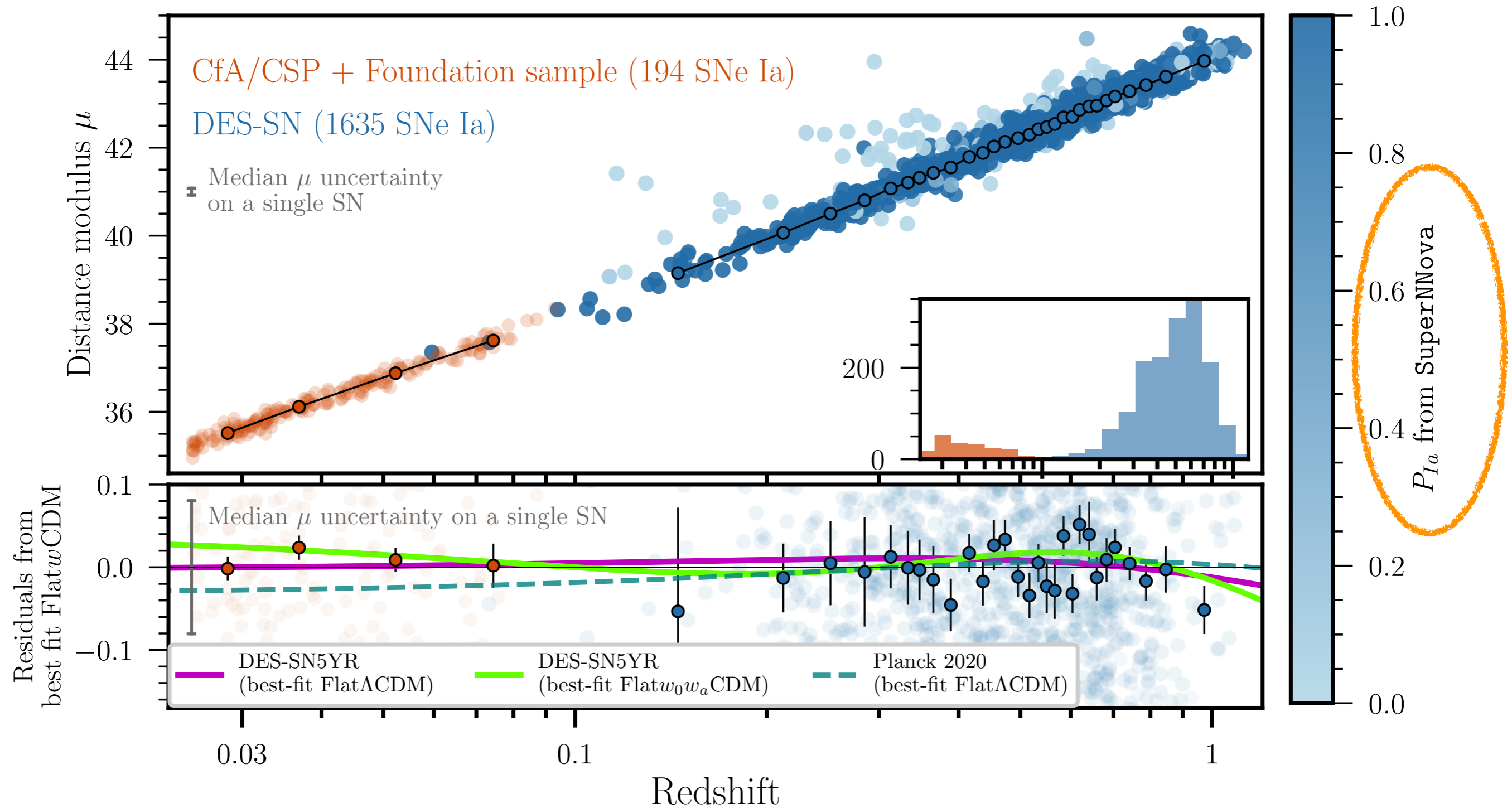
# DES: Light-curves + host-galaxy redshifts



**1484 SNe Ia**

*Largest high- $z$  SN Ia sample from a single survey for cosmology*

*Möller+ 2022a, Vincenzi+ 2024, DES+ 2024*

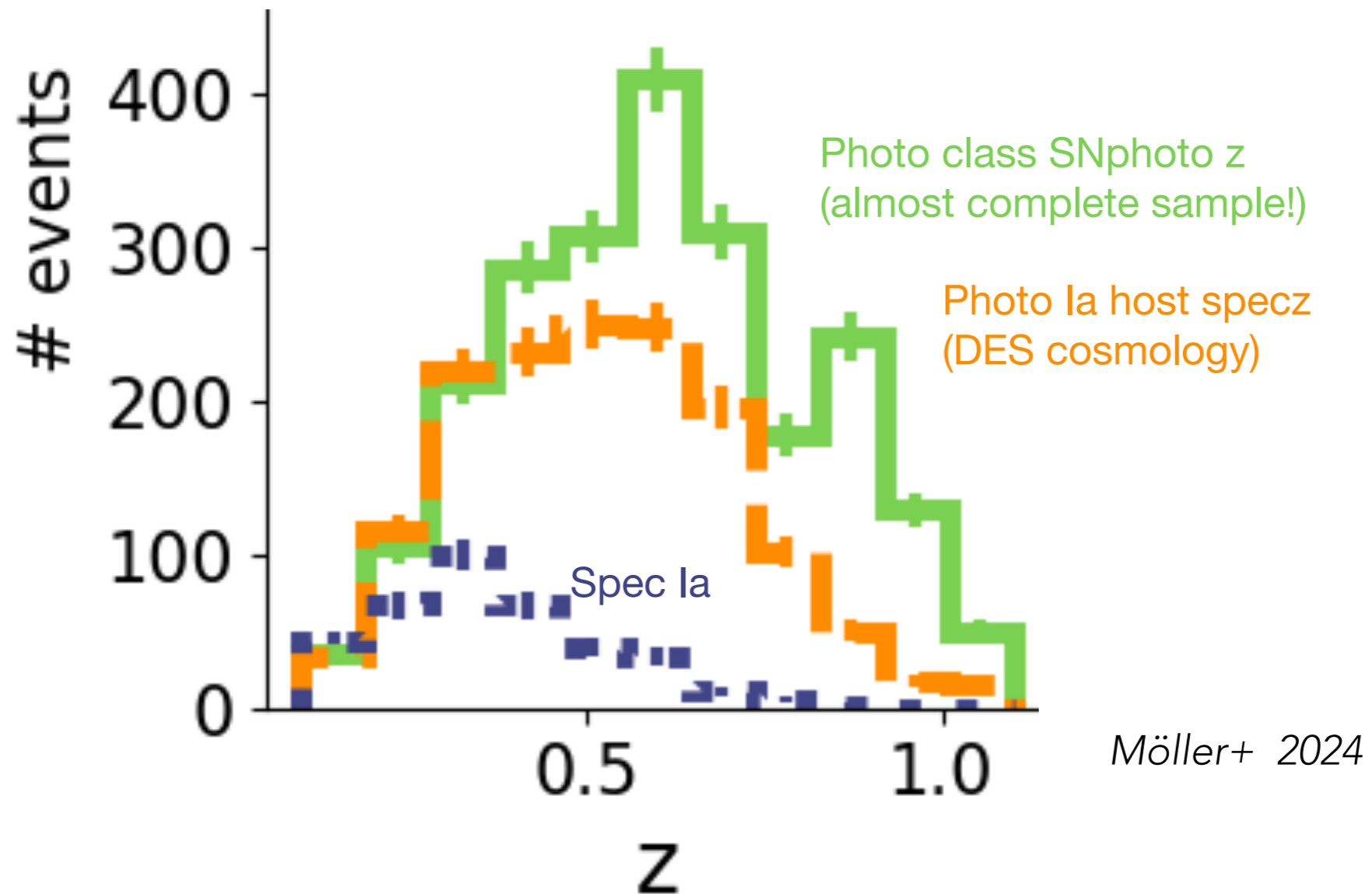


DES Collaboration 2024

# DES: Light-curves + ~~host-galaxy redshifts~~



# DES: Light-curves + ~~host-galaxy redshifts~~



**2,411 high quality SNe Ia**

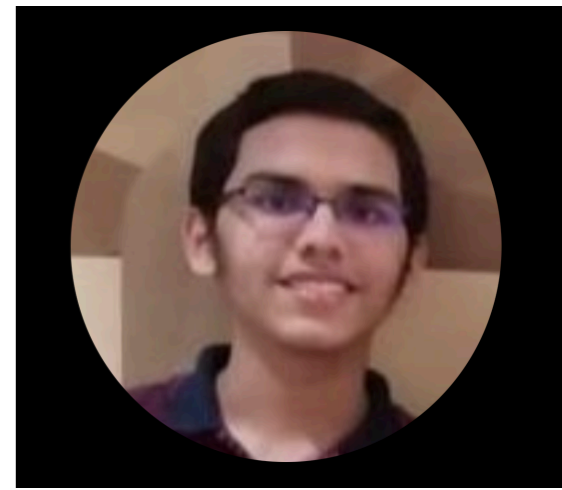
1. Supervised + simulations

**2. Unsupervised**

3. AL for improving training sets

# Unsupervised SN vs AGNs

Preliminary!!!



*Gondhalekar+ in prep.*

# <https://fink-portal.org/download>

## Fink Data Transfer

### Select data source

Source: ZTF

### Filter alerts

Dates: 2021-06-01 - 2022-06-01

Classe(s): ['(SIMBAD) AGN', '(SIMBAD) Blazar', '(SIMBAD) BLLac', '(SIMBAD) LINER', '(TNS) QSO', '(SIMBAD) QSO', '(SIMBAD) Seyfert', '(SIMBAD) Seyfert\_1', '(SIMBAD) Seyfert\_2', '(TNS) SN', '(TNS) SN I', '(TNS) SN Ia', '(TNS) SN Ia-91bg-like', '(TNS) SN Ia-91T-like', '(TNS) SN Ia-CSM', '(TNS) SN Ia-pec', '(TNS) SN Iax[02cx-like]', '(TNS) SN Ib', '(TNS) SN Ib-Ca-rich', '(TNS) SN Ib-pec', '(TNS) SN Ib/c', '(TNS) SN Ibn', '(TNS) SN Ic', '(TNS) SN Ic-BL', '(TNS) SN Ic-pec', '(TNS) SN Icn', '(TNS) SN II', '(TNS) SN II-pec', '(TNS) SN Iib', '(TNS) SN IIL', '(TNS) SN IIn', '(TNS) SN IIn-pec', '(TNS) SN IIP', 'Early SN Ia candidate']

Conditions: candidate.rb >= 0.65;  
candidate.nbad = 0; candidate.fwhm <= 5;  
candidate.elong <= 1.2;  
abs(candidate.magdiff) <= 0.1;

### Select content

Content: Lightcurve

### Submit

Trigger your job!

Description

Log in

Data Source

Choose the type of alerts you want to retrieve

ZTF  ELASTICC (v1)  ELASTICC (v2.0)  ELASTICC (v2.1)

Filters

Date Range \*

Pick up start and stop dates (included).

June 1, 2021 – June 1, 2022

Alert class

Select all classes you like! Default is all classes.

(SIMBAD) AGN x (SIMBAD) Blazar x (SIMBAD) BLLac x (SIMBAD) LINER x (TNS) QSO x (SIMBAD) QSO x (SIMBAD) Seyfert x (SIMBAD) Seyfert\_1 x (SIMBAD) Seyfert\_2 x (TNS) SN x (TNS) SN I x (TNS) SN Ia x (TNS) SN Ia-91bg-like x (TNS) SN Ia-91T-like x (TNS) SN Ia-CSM x (TNS) SN Ia-pec x (TNS) SN Iax[02cx-like] x (TNS) SN Ib x (TNS) SN Ib-Ca-rich x (TNS) SN Ib-pec x (TNS) SN Ib/c x (TNS) SN Ibn x (TNS) SN Ic x (TNS) SN Ic-BL x (TNS) SN Ic-pec x (TNS) SN Icn x (TNS) SN II x (TNS) SN II-pec x (TNS) SN Iib x (TNS) SN IIL x (TNS) SN IIn x (TNS) SN IIn-pec x (TNS) SN IIP x (Fink) Early Supernova Ia candidates x

Extra conditions

One condition per line (SQL syntax), ending with semi-colon. See [here](#) (and also [here](#)) for fields description and [here](#) for examples.

```
candidate.rb >= 0.65;  
candidate.nbad = 0;  
candidate.fwhm <= 5;  
candidate.elong <= 1.2;  
abs(candidate.magdiff) <= 0.1;
```



Recommended filters from the ZTF team to yield a fairly pure sample

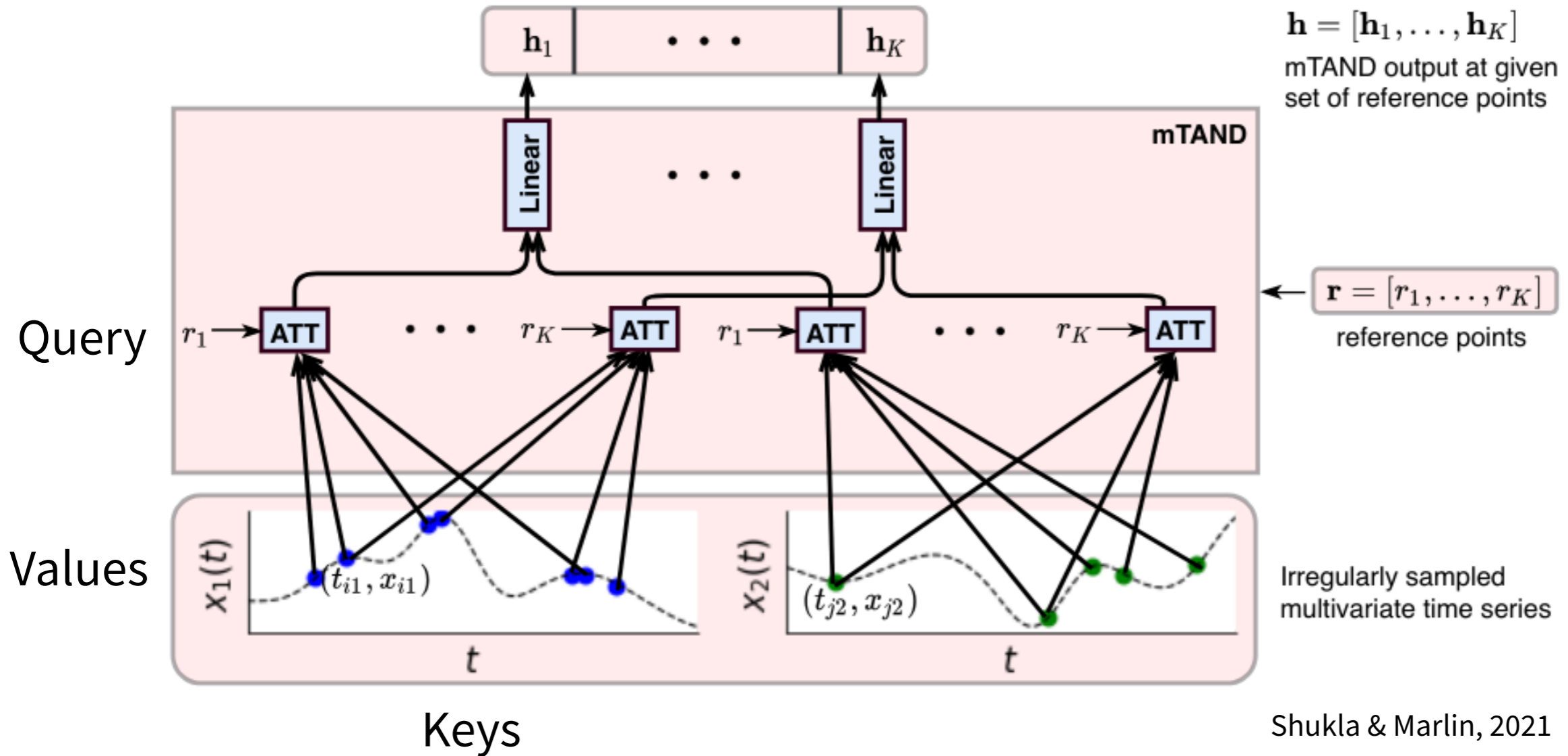
Alert content

Choose the content you want to retrieve

Lightcurve (~1.4 KB/alert)  Cutouts (~41 KB/alert)  Full packet (~55 KB/alert)

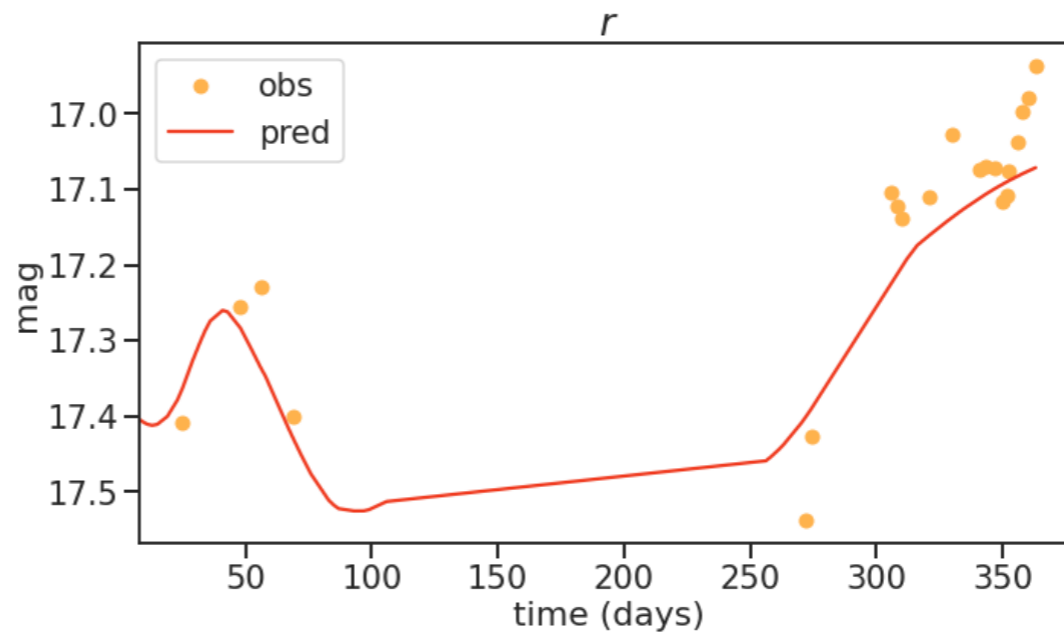
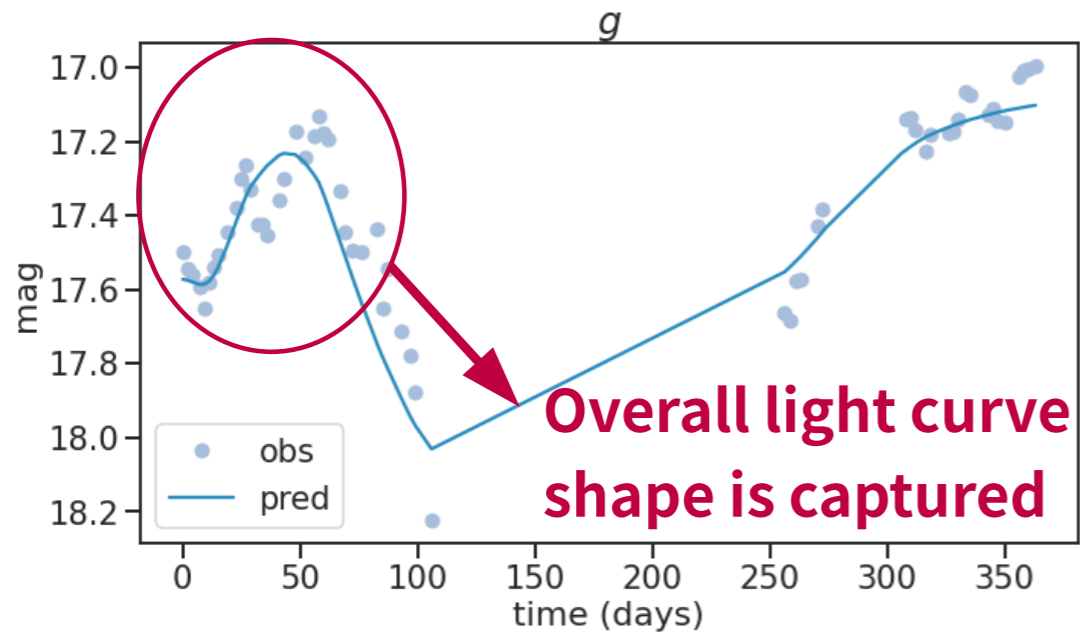
Submit job

# mTAN: Multi-Time Attention Network

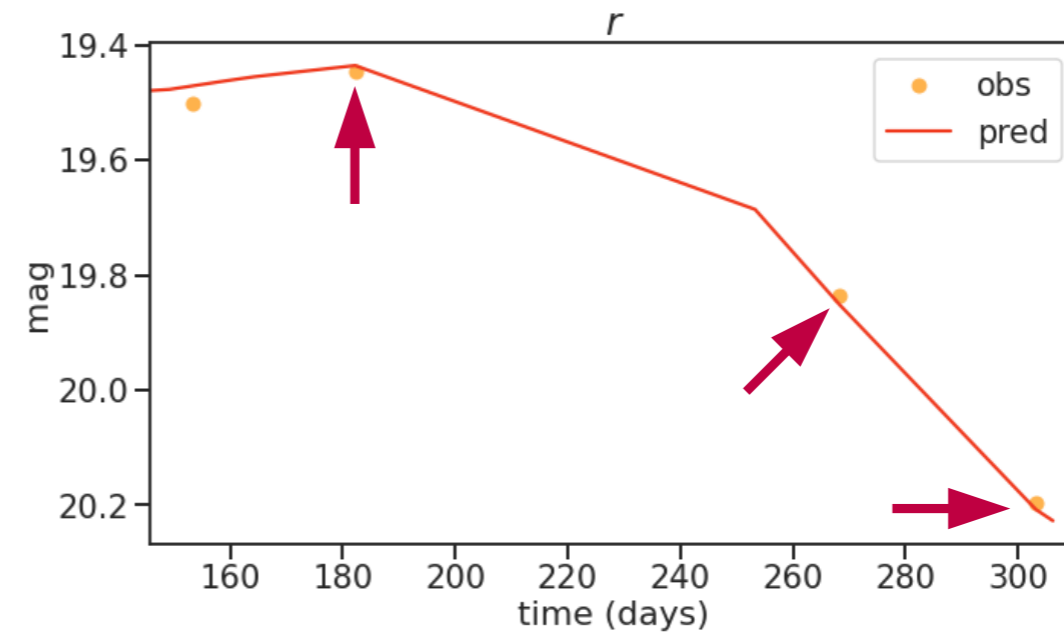
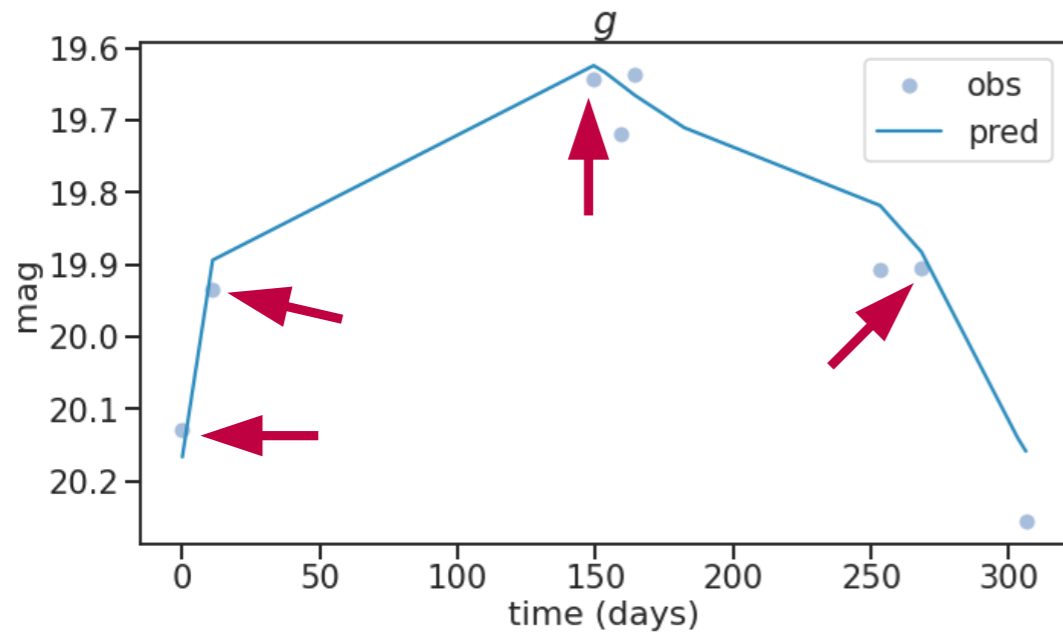


**Preliminary!!!**

ZTF18acaajdo - Seyfert\_1

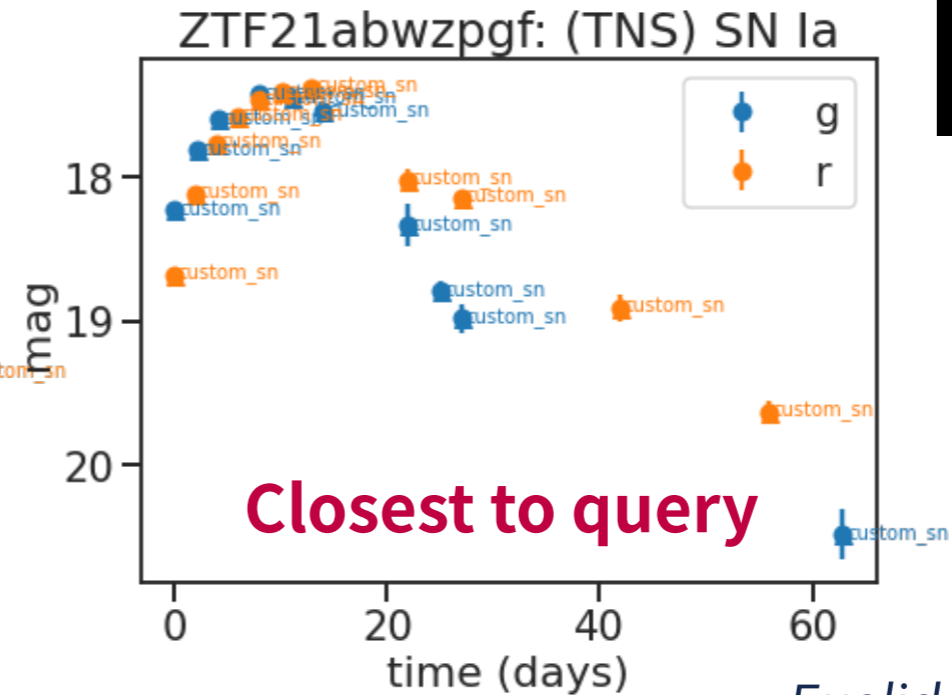
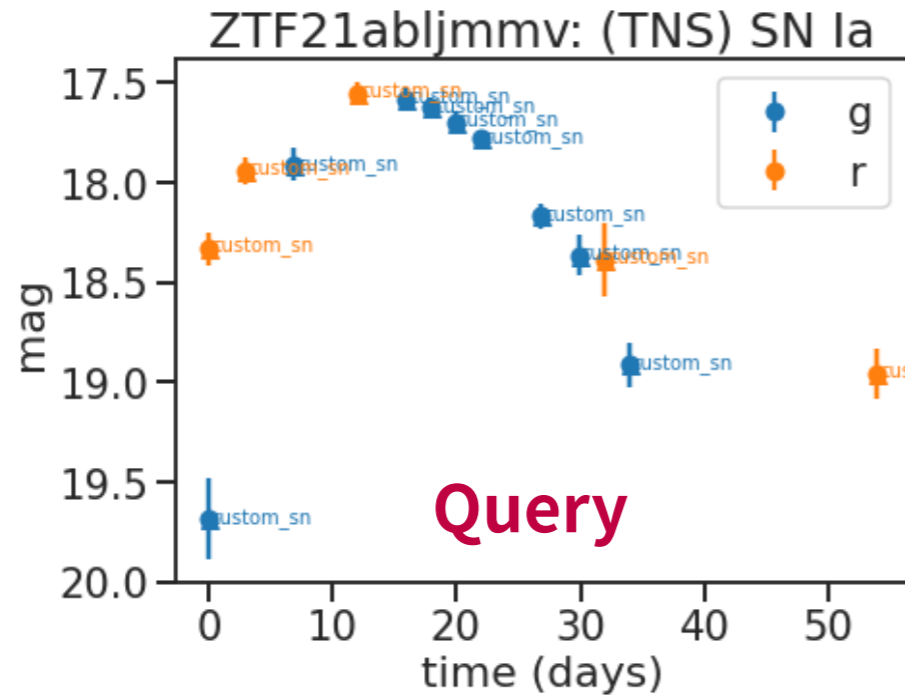
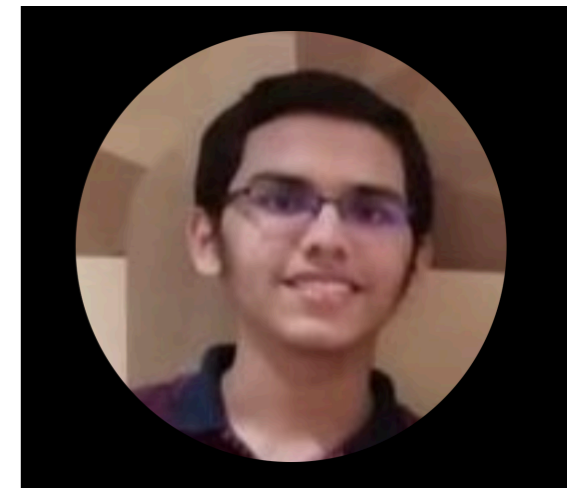


ZTF18aakonwr - Unknown

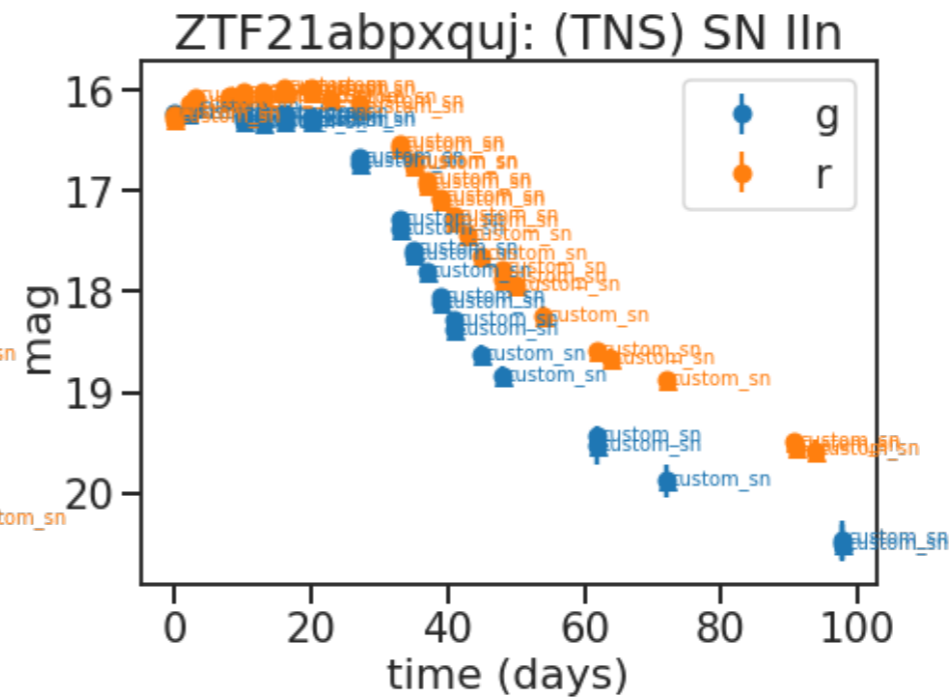
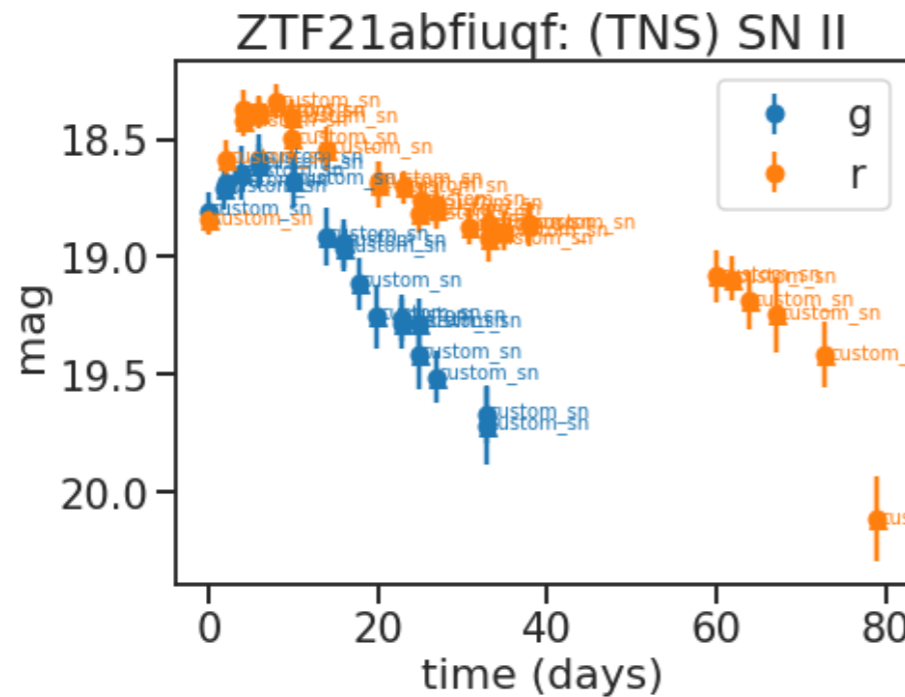


*Gondhalekar+ in prep.*

Preliminary!!!

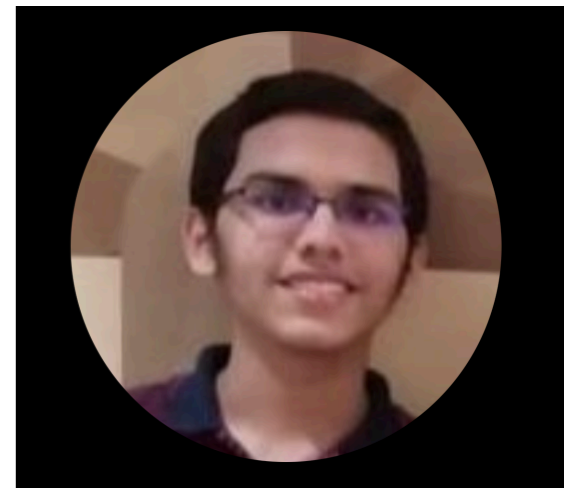


*Euclidean, cosine distance...*



*Gondhalekar+ in prep.*

**~0.019 s /light-curve**



Compared to the time-series transformer model of Allam, Jr. & McEwen (2021), this is\*:

1. ~80 times faster inference (0.019s vs. 1.5s).
2. ~3 times smaller model size (340 KB vs. 1100 KB).

Assuming ~1 million light curves/year, this amounts to ~0.34 kg of CO<sub>2</sub>, equivalent to driving 1.37 km by an average ICE car\*\*.

\*The follow-up paper (Allam, Jr. et al 2023) reduced model size by ~18 and execution time by ~8 times than their original approach.

\*\*Source: [MachineLearning Impact calculator](#)

*Gondhalekar+ in prep.*



1. Supervised + simulations

2. Unsupervised

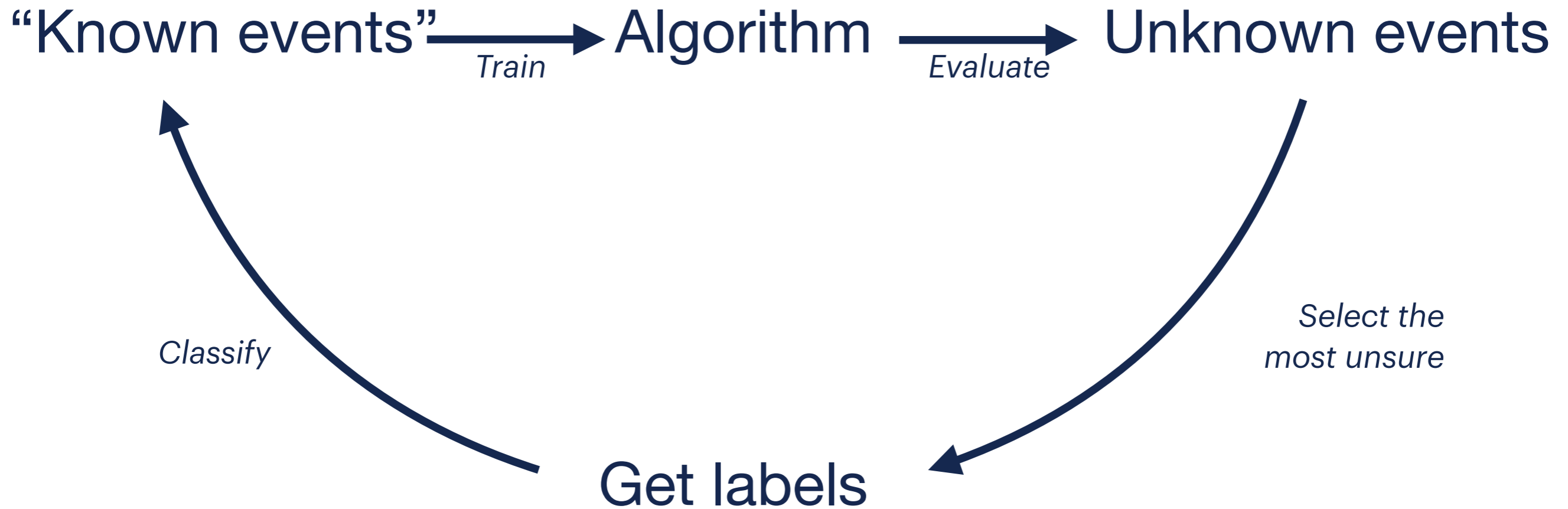
**3. AL for improving training sets**

“Known events”  $\xrightarrow{\text{Train}}$  Algorithm  $\xrightarrow{\text{Evaluate}}$  Unknown events



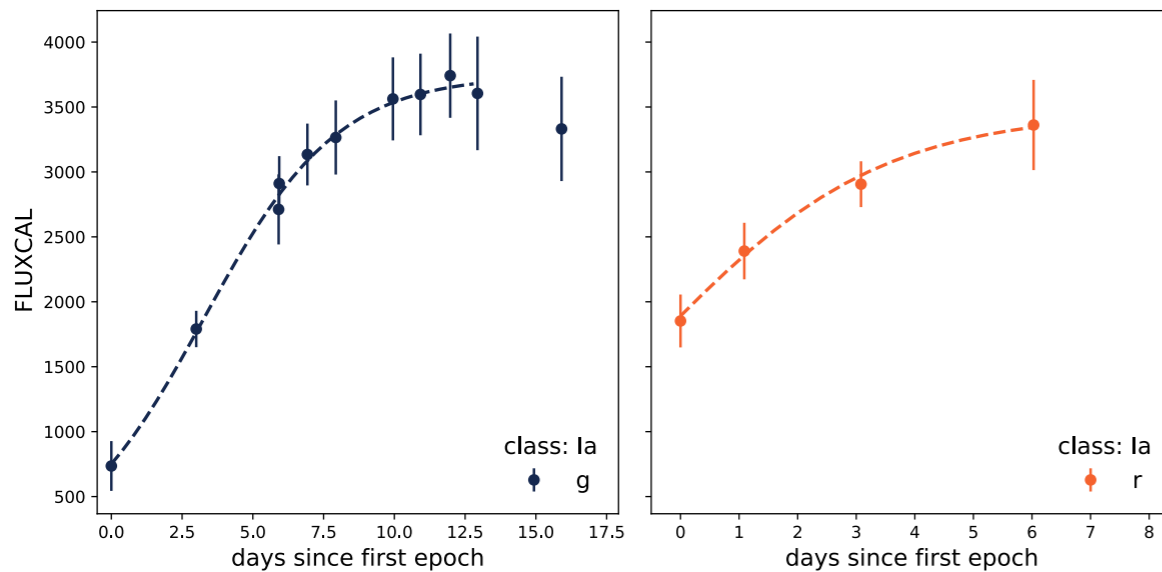
Your target  
(e.g. SNIa)

# Active Learning



# Early SNe Ia

We don't know SNe Ia progenitors  
Missing modelling for early SNe Ia

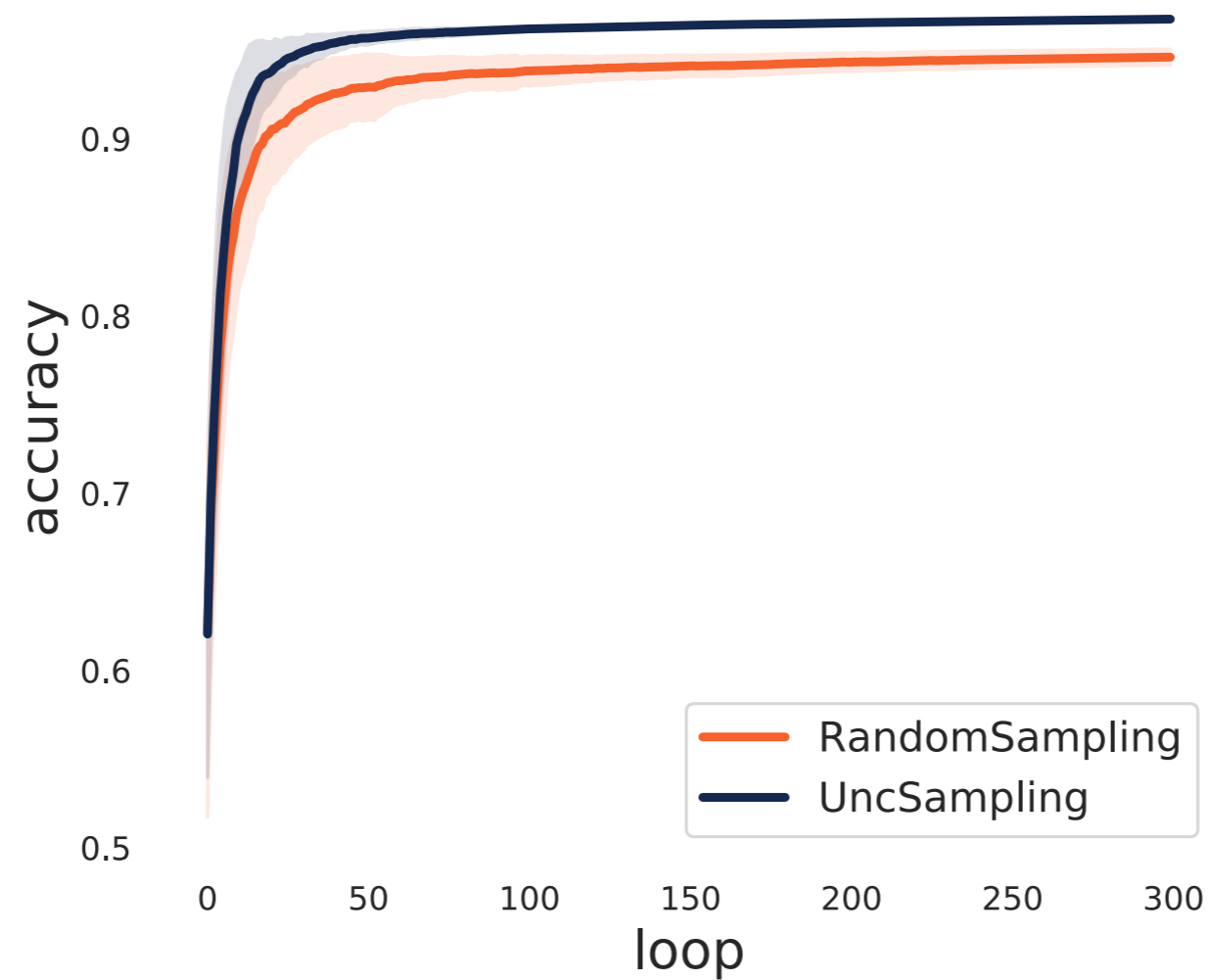


Features: sigmoid fit  
Random Forest

*Ishida+2019*

*Leoni, Ishida+ 2022*

# Active Learning SN classification (simulated)



*Leoni, Ishida et al. 2022*

# Active Learning SN classification (real data)

~30 known Ics

*Leoni, Ishida et al. 2022*

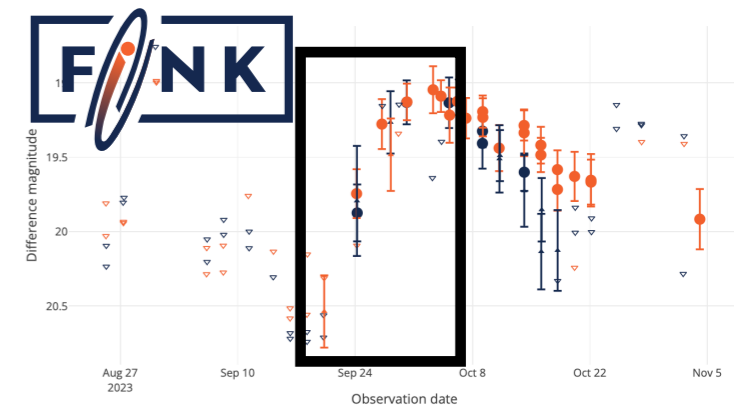


*Train*

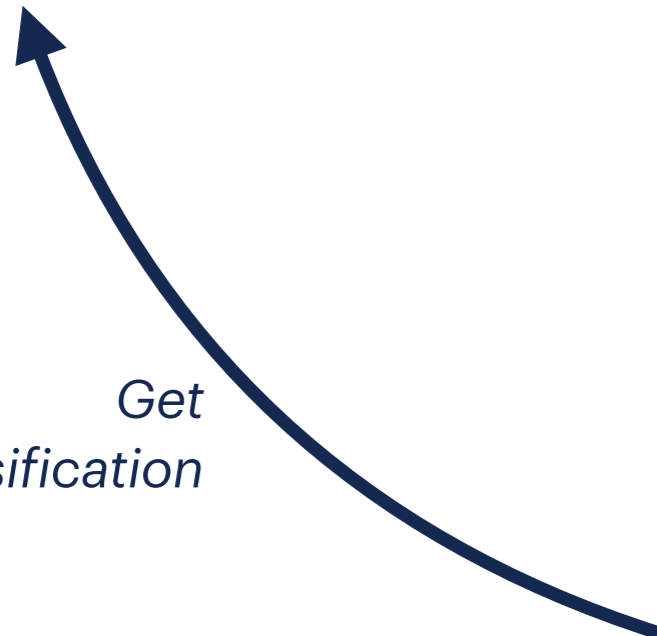
RF



*Evaluate*

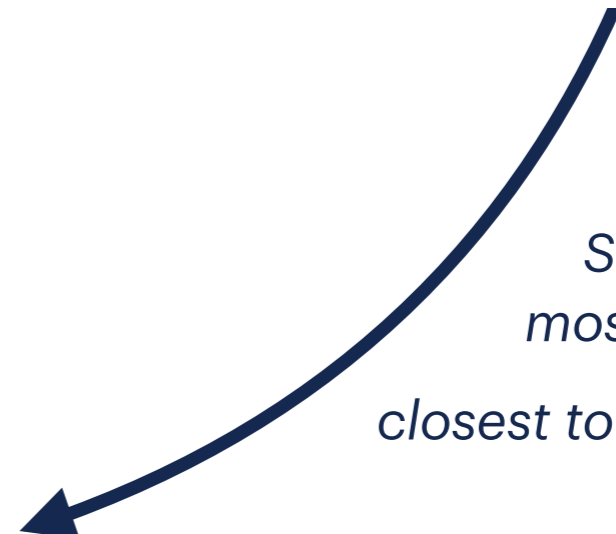


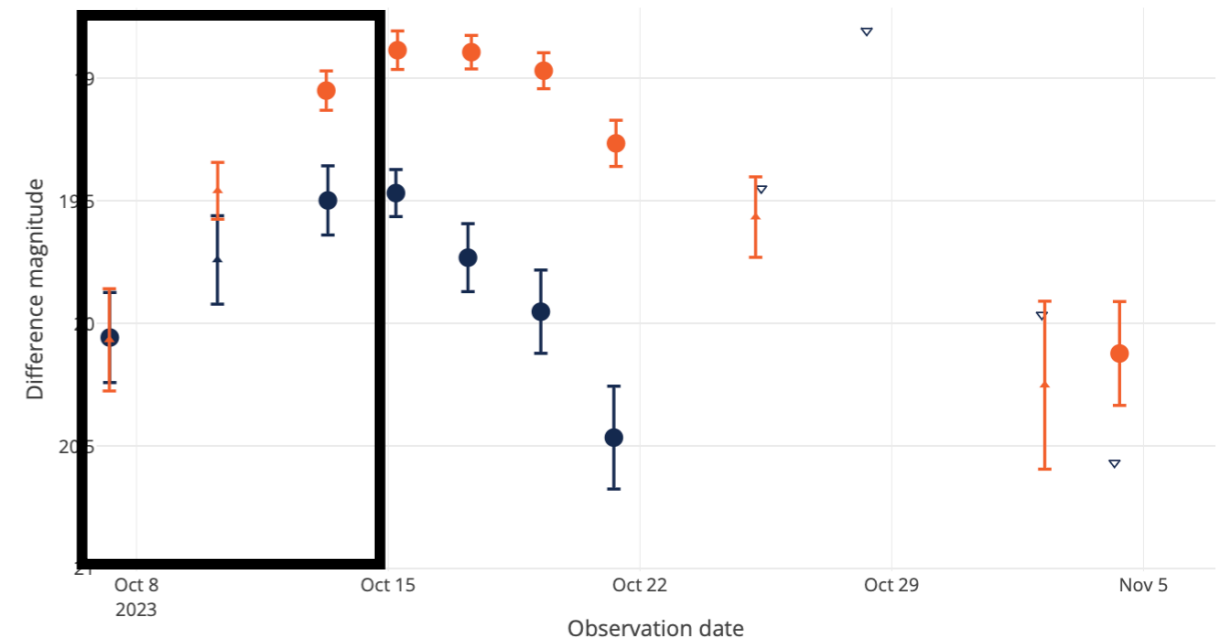
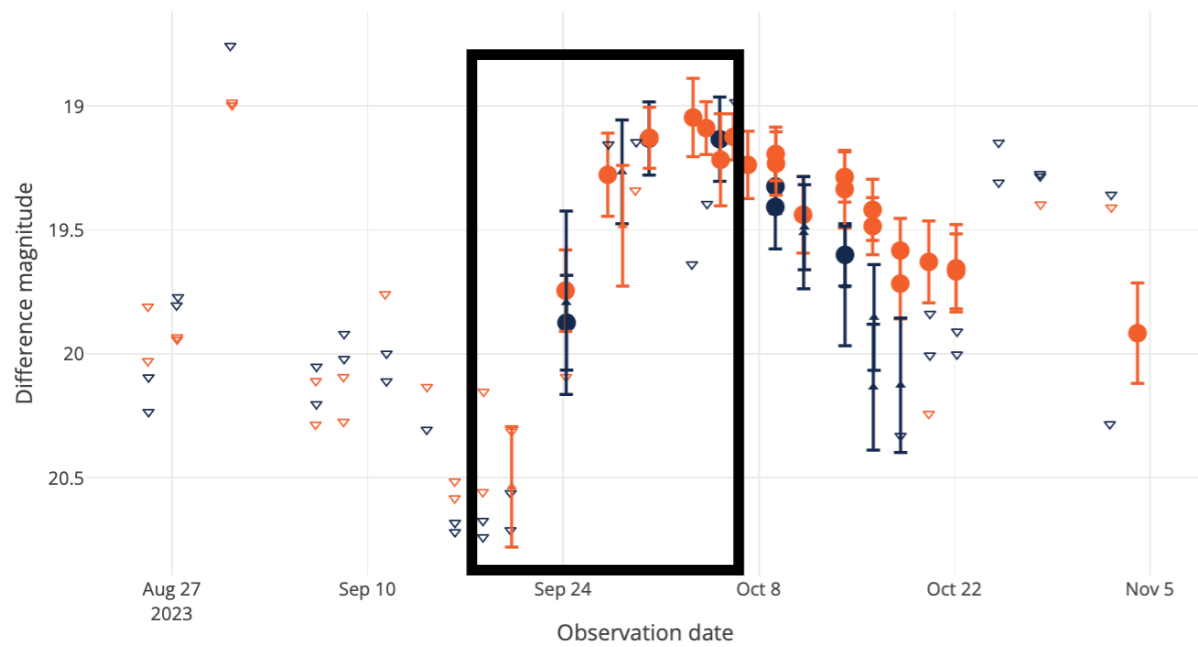
*Get  
classification*



Spectroscopy

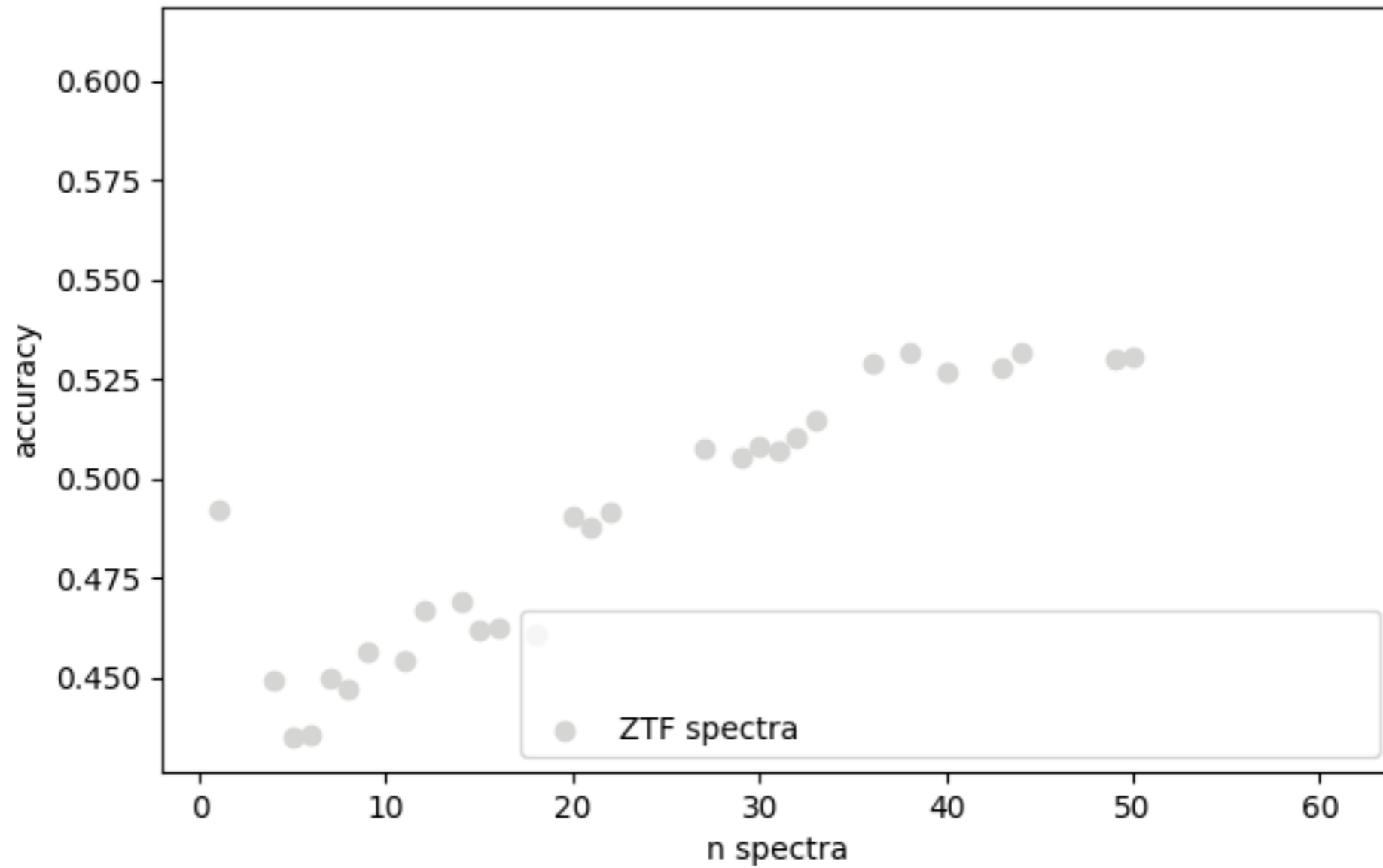
*Select the  
most unsure  
closest to  $P_{Ia} = 0.5$*





2 CV, 16 SN Ia, 5 SN II, 1 SN Ib, 2 SN Ibn, 1 SN Ic, 1 microlensing event...

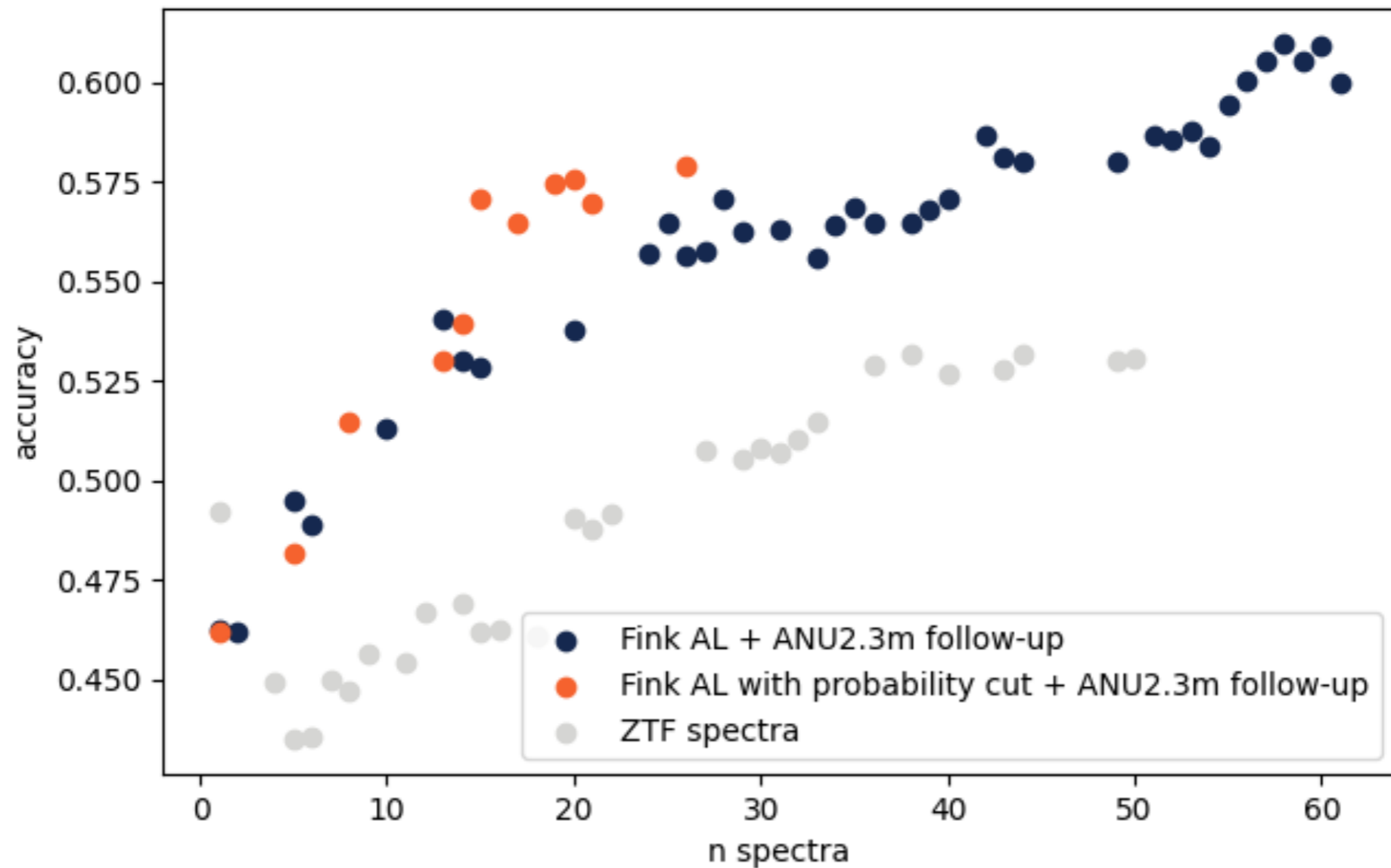
Preliminary!!!



*40 light-curves as  
initial training set*



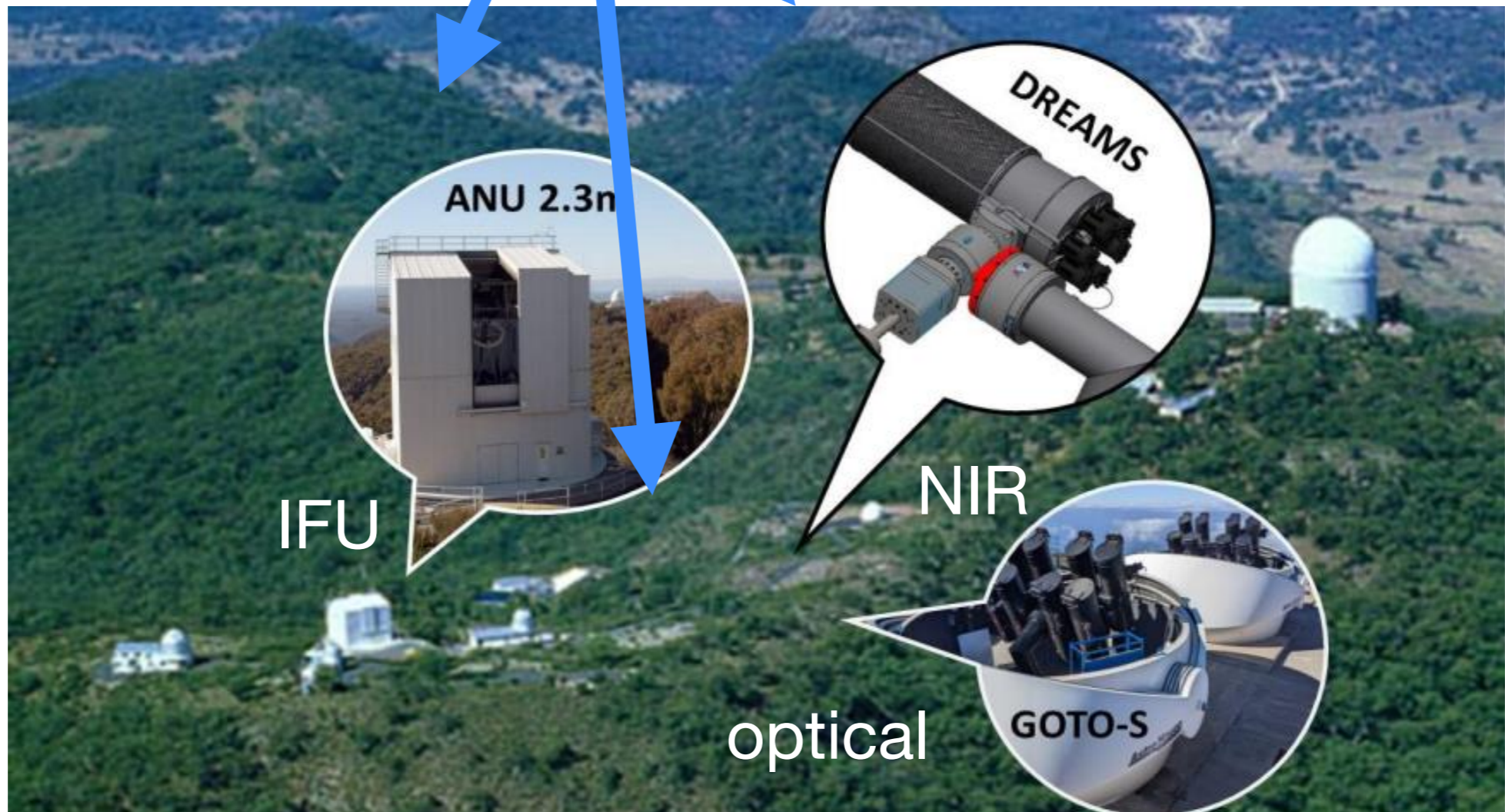
Preliminary!!!



61 labels  
26 labels  
50 labels

# Robotic Network @ Siding Spring Observatory

Lidman, Colless, Wolf, Travouillon, Brough, Seitzzahl, Ruiter, Brown, Heger, Galloway, Möller, Kamath, O'Toole, Wen, Einecke



Automatic data reduction

*Rubin is an amazing opportunity and a big data challenge!*



- All our data is public [fink-portal.org](https://fink-portal.org)
- First candidate ID modules: SNe, KNe, SSO, GRB, microlensing...

### *Supernova ML challenges:*

- AL is a good strategy to improve training sets (optimise follow-up)
- This will percolate to simulations (built with templates)
- And then we can use large NNs accurately for precision science!



# Backup

